# A Point-Based Algorithm for Multiple 3D Surface Alignment of Drug-Sized Molecules

**Dissertation**

am Fachbereich Mathematik und Informatik
der Freien Universität Berlin

zur Erlangung des akademischen Grades
des Doktors der Naturwissenschaften (Dr. rer. nat.)
vorgelegt von

**Diplom-Informatiker**

**Daniel Baum**

Berlin, März 2007

# Abstract

One crucial step in *virtual drug design* is the identification of new *lead structures* with respect to a pharmacological target molecule. The search for new lead structures is often done with the help of a *pharmacophore*, which carries the essential structural as well as physico-chemical properties that a molecule needs to have in order to bind to the target molecule. In the absence of the target molecule, such a pharmacophore can be established by comparison of a set of active compounds. In order to identify their common features, a *multiple alignment* of all or most of the active compounds is necessary. Moreover, since the "outer shape" of the molecules plays a major role in the interaction between drug and target, an alignment algorithm aiming at the identification of common binding properties needs to consider the molecule's "outer shape", which can be approximated by the *solvent excluded surface*.

In this thesis, we present a new approach to molecular surface alignment based on a discrete representation of shape as well as physico-chemical properties by points distributed on the solvent excluded surface. We propose a new method to distribute points regularly on a surface w.r.t. a smoothly varying point density given on that surface. Since the point distribution algorithm is not restricted to molecular surfaces, it might also be of interest for other applications. For the computation of pairwise surface alignments, we extend an existing *point matching scheme* to surface points, and we develop an efficient data structure speeding up the computation by a factor of three. Moreover, we present an approach to compute multiple alignments from pairwise alignments, which is able to handle a large number of surface points. All algorithms are evaluated on two sets of molecules: eight *thermolysin inhibitors* and seven *HIV-1 protease inhibitors*. Finally, we compare the results obtained from surface alignment with the results obtained by applying an atom alignment approach.

# Zusammenfassung

Die Identifizierung neuer *Leitstrukturen* (lead structures) zur Entwicklung optimierter Wirkstoffe ist ein äußerst wichtiger Schritt in der *virtuellen Wirkstoffentwicklung* (virtual drug design). Die Suche nach neuen Leitstrukturen wird oft mit Hilfe eines *Pharmakophor-Modells* durchgeführt, welches die wichtigsten strukturellen wie auch physiko-chemischen Eigenschaften eines bindenden Moleküls in sich vereint. Ist das Zielmolekül (target) nicht bekannt, kann das Pharmakophor-Modell mit Hilfe des Vergleiches aktiver Moleküle erstellt werden. Hier ist insbesondere die *gleichzeitige Überlagerung* (multiple alignment) aller oder nahezu aller Moleküle notwendig. Da bei der Interaktion zweier Moleküle die "äußere Form" der Moleküle eine besondere Rolle spielt, sollte diese von jedem Überlagerungsalgorithmus, der sich mit der Identifizierung von Bindungseigenschaften befasst, berücksichtigt werden. Dabei kann die "äußere Form" durch eine bestimmte Art von molekularer Oberfläche approximiert werden, die man als *solvent excluded surface* bezeichnet.

In dieser Arbeit stellen wir einen neuen Ansatz zur Überlagerung molekularer Oberflächen dar, der auf einer diskreten Repräsentation sowohl der Form als auch der molekularen Eigenschaften mittels Punkten beruht. Um die Punkte auf der molekularen Oberfläche möglichst regulär entsprechend einer gegebenen Punktdichte zu verteilen, entwickeln wir eine neue Methode. Diese Methode ist nicht auf Moleküloberflächen beschränkt und könnte daher auch für andere Anwendungen von Interesse sein. Basierend auf einem bekannten *Point-Matching* Verfahren entwickeln wir einen Point-Matching Algorithmus für Oberflächenpunkte. Dazu erarbeiten wir u.a. eine effiziente Datenstruktur, die den Algorithmus um einen Faktor von drei beschleunigt. Darüberhinaus stellen wir einen Ansatz vor, der Mehrfachüberlagerungen (multiple alignments) aus paarweisen Überlagerungen berechnet. Die Herausforderung besteht hierbei vor allem in der großen Anzahl von Punkten, die berücksichtigt werden muss. Die vorgestellten Algorithmen werden an zwei Gruppen von Molekülen evaluiert, wobei die erste Gruppe aus acht *Thermolysin Inhibitoren* besteht, die zweite aus sieben *HIV-1 Protease Inhibitoren*. Darüberhinaus vergleichen wir die Ergebnisse der Oberflächenüberlagerung mit denen einer Atommittelpunktüberlagerung.

# Acknowledgments

# Contents